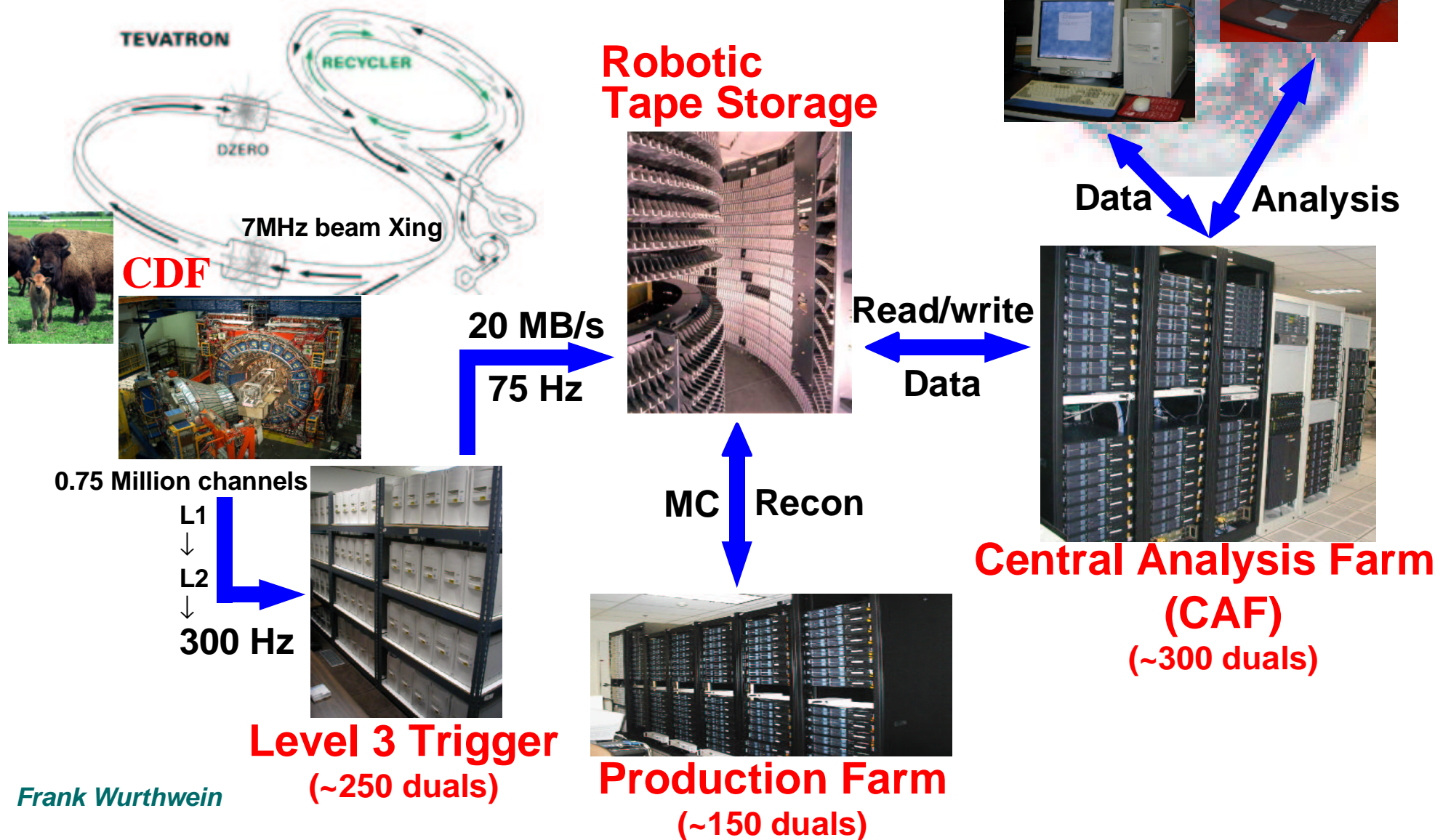# User Analysis Computing at CDF

## Frank Wurthwein
*MIT/UCSD/FNAL-CD*
for the CDF Collaboration

➢ **Computing Model**

   ➢ **status**

   ➢ **Future directions**

# CDF DAQ/Analysis Flow

**User Desktops**

**TEVATRON**

**RECYCLER**

**DZERO**

**7MHz beam Xing**

**CDF**

**Robotic Tape Storage**

**Data** **Analysis**

**20 MB/s**

**75 Hz**

**Read/write**

**Data**

**0.75 Million channels**

L1
↓
L2
↓
**300 Hz**

**MC** **Recon**

**Central Analysis Farm (CAF)**
**(~300 duals)**

**Level 3 Trigger**
**(~250 duals)**

**Production Farm**
**(~150 duals)**

*Frank Wurthwein*

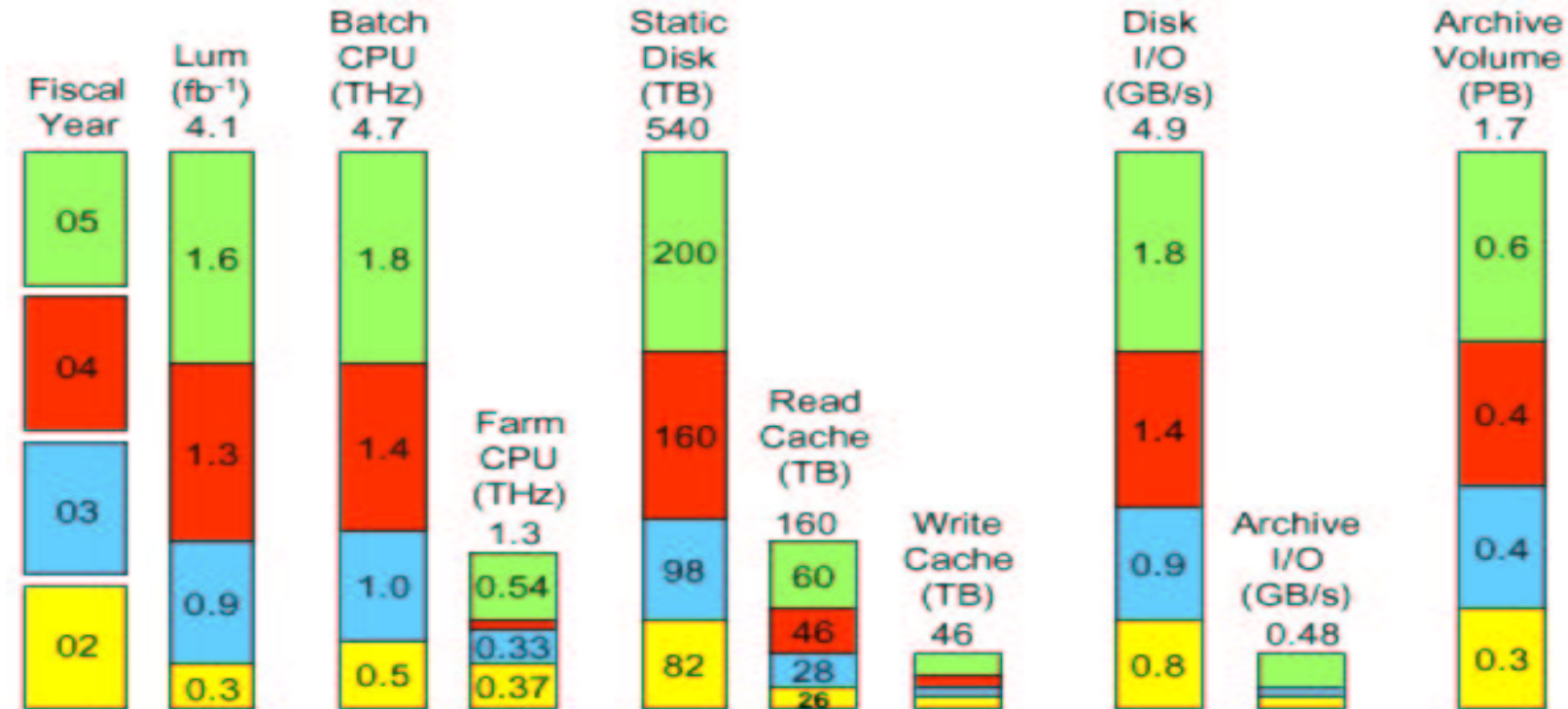# **Data/Software Characteristics**

## **Data Characteristics:**

- **Root I/O: ~80-400 kB/event (configurable content)**
- **'Standard' ntuple: 5-10 kB/event**
- **Typical RunIIa secondary dataset size: $10^7$ events**
- **Winter03 physics: ~100 datasets adding up to ~50TB**
- **Largest dataset for Winter03 physics: 3.5e7 evts**

## **Analysis Software:**

- **Typical analysis jobs run @ 5 Hz on 1 GHz P3**
    $\rightarrow$ **few MB/sec**
- **CPU rather than I/O bound (FastEthernet)**

# Computing Requirements



**Requirements set by goal:**

200 simultaneous users to analyze secondary data set ($10^7$ evts) in a day

**Need ~700 TB of disk and ~5 THz of CPU by end of FY'05:**

$\rightarrow$ **need lots of disk** $\rightarrow$ **need cheap disk** $\rightarrow$ **IDE Raid**

$\rightarrow$ **need lots of CPU** $\rightarrow$ **commodity CPU** $\rightarrow$ **dual Intel/AMD**

# Computing Model

**Interactive Computing on desktop:**
- ➢ **Complete access to all data from desktop via dCache & rootd**

**Batch Computing on "remote" cluster(s):**
- ➢ **Binary compatible with desktop**
- ➢ **qsub, qstat, kill, ls, tail, top via command line/web**
- ➢ **Large scale parallelisation with single submission**
  - ➔ **Single summary email upon completion**
- ➢ **User scratch space inside cluster**
  - ➔ **Krb5 ticket created @ launch time**
- ➢ **Data access Winter03: 90% NFS+rootd, 10% dCache**

# Example job submission

- Compile, build, debug analysis job on 'desktop'

- Fill in appropriate fields & submit job

- Retrieve output using kerberized FTP tools ... or write output directly to 'desktop'!

section integer range

CDF RunII CAF GUI

Initial Command: ./simple.sh$    600  610

Process Type: Short

Original Directory: /home/msn/releases/development/CafUtil/examples   Browse...

Ouput File Location: msn@fcdflnx2.fnal.gov/cdf/scratch/msn/temp$.tgz

Email?   Email Address: msn@fnal.gov

Submit   Quit                                      Ready

(2002-05-23 01:46:51) Email sent to msn@fnal.gov upon job completion
(2002-05-23 01:46:55) /bin/tar -cvzf /home/msn/msn49959.tgz *
(2002-05-23 01:46:57) Remove /home/msn/msn49959.tgz
(2002-05-23 01:46:57) Job Submission is successful, JID: 873

output destination          user exe+tcl directory

# Web Monitoring of User Queues

**Each user a different queue**

**Process type for job length**
- **test**: 5 mins
- **short**: 2 hrs
- **medium**: 6 hrs
- **long**: 2 days

**This example:**
1 job → 11 sections

**(+ 1 additional section automatic for job cleanup)**



Netscape: FBSWWW CAF list of queues

File   Edit   View   Go   Communicator                                          Help

**FBSNG on the web**
Farm:          CAF
Time:          Thu May 23 02:32:41 2002
Report:        List of queues

Queues  Jobs  Nodes  Process Types

User Monitor

| Name | Status | Default Process Type | Share | Prio | Waiting | Ready | Running | Total |
|------|--------|----------------------|-------|------|---------|-------|---------|-------|
| akom | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| amitl | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| anikeev | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| belforte | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msmartin | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msn | OK | short | 1.00 | 0 | 1 | 0 | 11 | 12 |
| pauly | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| paus | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| ratnikov | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| rescigno | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| semenia | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sfiligoi | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sgromoll | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| shepard | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sidoti | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| spezziga | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| test | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thkim | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thom | OK | short | 1.00 | 0 | 1 | 0 | 1 | 2 |

# Monitoring jobs in your queue

# Monitoring sections of your job

CAF Hardware Architecture
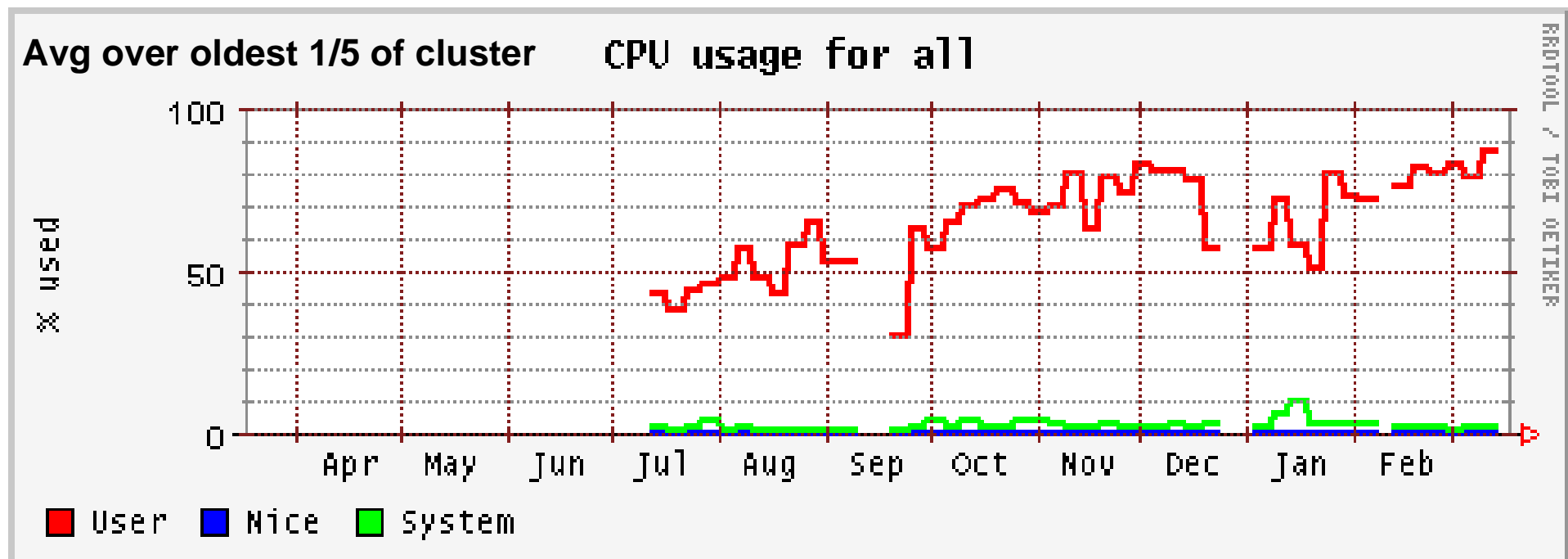
# CAF utilization

## User perspective:
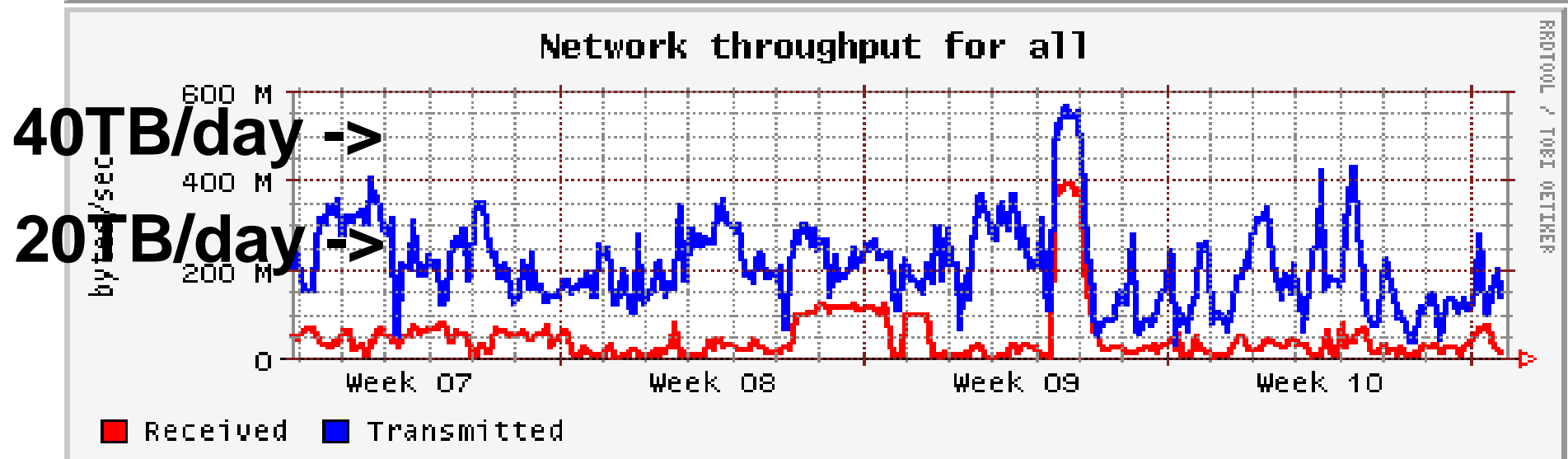- **10,000 jobs launched/day**
- **400 users total**
- **100 users per day**

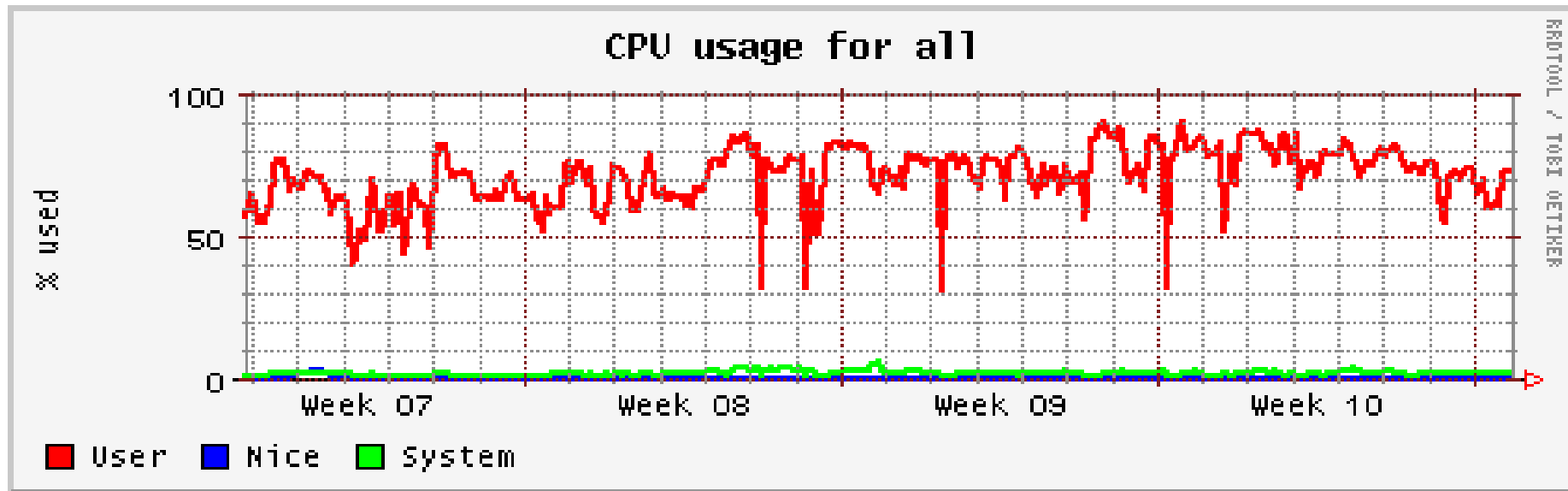## System perspective:
- **Up to 90% avg CPU utilization**
- **200-600MB/sec I/O**
- **Failure rate ~1/2000**
- **Avg uptime of WN = 60days**



Avg over oldest 1/5 of cluster — CPU usage for all

Legend: ■ User ■ Nice ■ System

# CAF utilization last month



**CPU usage for all**

**Network throughput for all**

40TB/day ->

20TB/day ->

# Status @ FNAL Today

User analysis computing based on commodity PC's

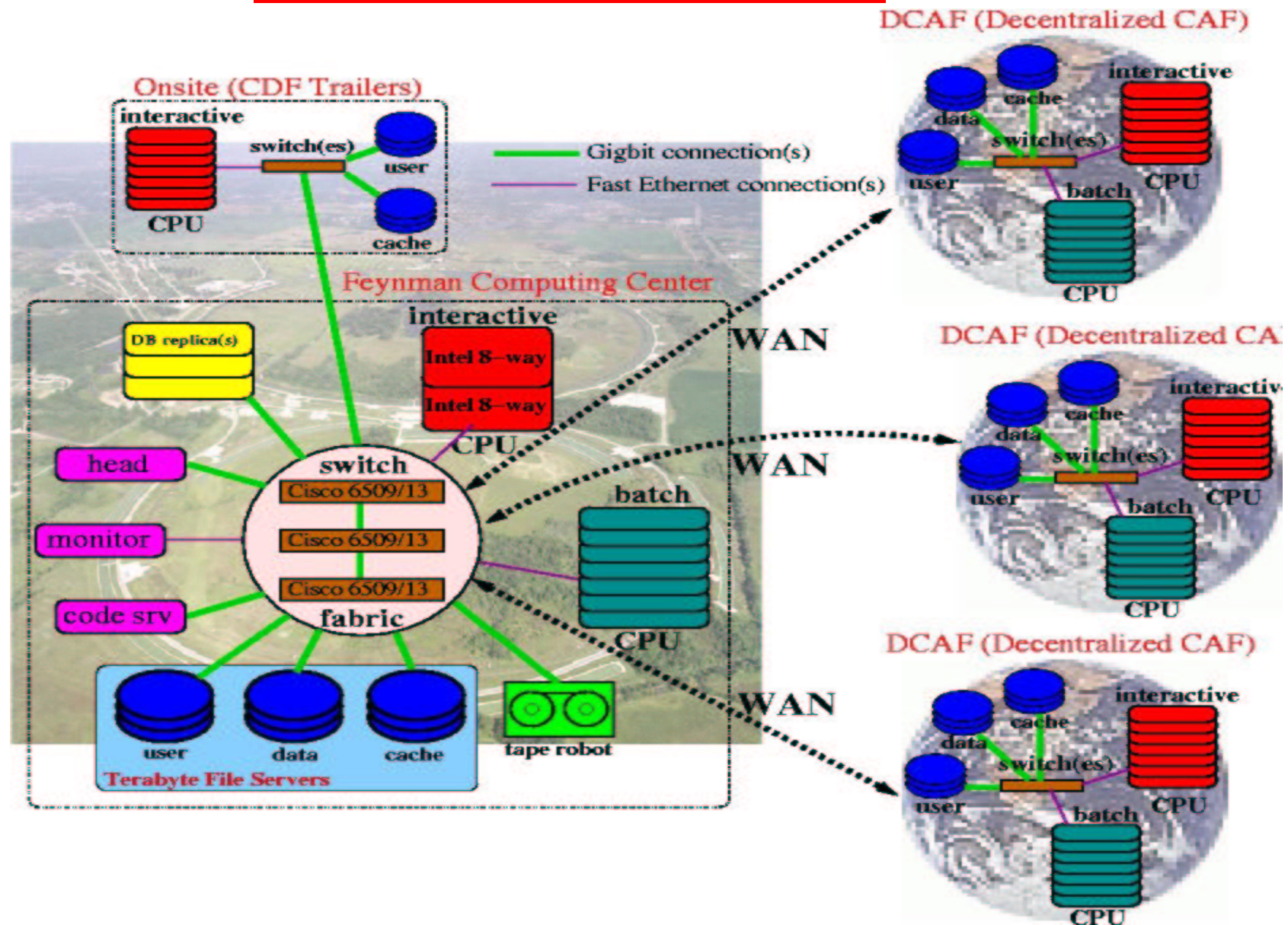**180TB disk space**      **1THz batch CPU**

Focus on building strong infrastructure

**up to 600MB/sec I/O**    **99.95% reliability**

that has been deployed as part of CDF grid "proof of principle" for SC2002 demo.

# Future Directions

# CDF grid = 3 pieces

**CAF:**
- Local cluster management
- Remote submission
- Fully in production: 99.95% reliability

**SAM:**
- WAN capable DH system
- Use for remote MC production in summer 03

**JIM:**
- Grid broker (based on condor,globus,sam)
- Proof of principle fall 02

# Summary & Conclusions

**CDF's computing model makes offsite computing contributions possible.**

**Offsite contributions, and accounting thereof is desirable.**

**The devil is in the detail.**